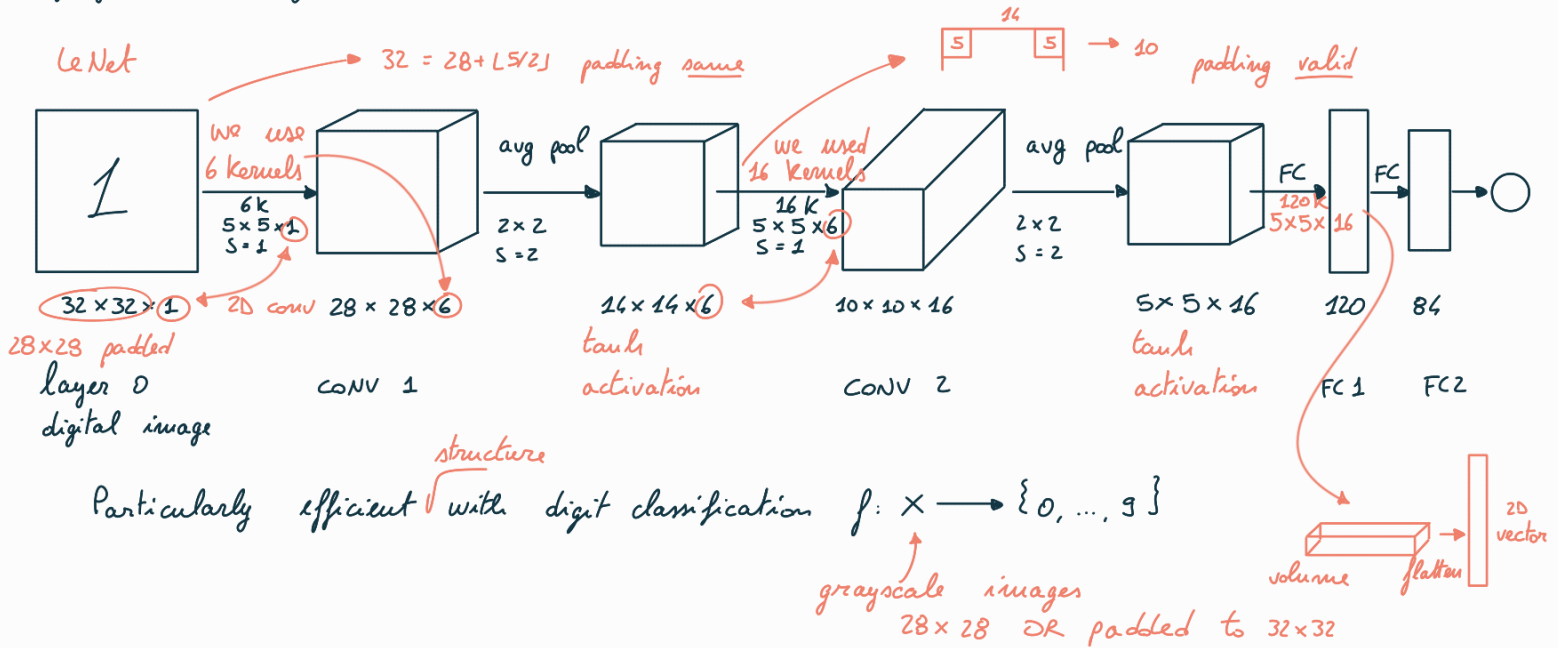


Convolutional Neural Networks for images (2D input)

Every convolutional layer of a CNN transforms the 3D input volume to a 3D output volume of neuron activations. We iterate this process several times and at the end the goal is to reach a situation in which we have a smaller volume. Usually the size of the volume after some levels is much smaller than the size of the input. Once we have this situation we can do what's called a **flatten operation**, that will produce just a vector of real values; we will continue with a standard fully connected layer.



2 convolutional layers + 2 fully connected layers

CONV 1	6k	5x5x1	→	6 · 5 · 5 + 6	bias
CONV 2	16k	5x5x6	→	16 · 5 · 5 · 6 + 16	
CONV 3	120k	5x5x16	→	120 · 5 · 5 · 16 + 120	
FC 1	120	84	→	120 · 84 + 84	
FC 2	84	10	→	84 · 10 + 10	

input output

trainable parameters

They are usually few at the beginning but they increase as the network goes deeper.

ILSVRC

2012	AlexNet	(8 layers)
2013	GoogleNet	(8 layers)
2014	VGG	(19 layers)
2015	ResNet	(22 layers)

Alexnet was the winner of ImageNet Large Scale Visual Recognition Competition (ILSVRC) 2012 (8 layers splitted in two parts trained in parallel to speed up computation - not required with modern hardware)

ResNet was the winner of ILSVRC in 2015. It introduced a new technique: **skip connections**. With skip connections you skip some of the layers and directly propagate the values some layers further without changes

